

Increasing Fault Resiliency in a Message-Passing Environment

Rolf Riesen, Kurt Ferreira

Ron Oldfield, Jon Stearley

James Laros, Kevin Pedretti

Ron Brightwell

rolf@sandia.gov

Sandia National Laboratories

Todd Kordenbrock

Hewlett-Packard Company

October 14, 2009

Motivation

Design

Evaluation

Analysis

Implications

Summary

Motivation

- Checkpoint/Restart is common way to deal with faults
- What can we do to increase checkpoint interval?
- Explore redundant computation for MPI applications
 - ◆ Can it be done at user level?
 - ◆ What are requirements for RAS system and runtime?
 - ◆ What is the overhead
 - ◆ Cost versus benefit?
- Write *rMPI* library at MPI profiling layer to learn what the issues are

Motivation

Design

Redundancy

Basics

Msg. order

Other

Status

Evaluation

Analysis

Implications

Summary

Design

Motivation

Design

Redundancy

Basics

Msg. order

Other

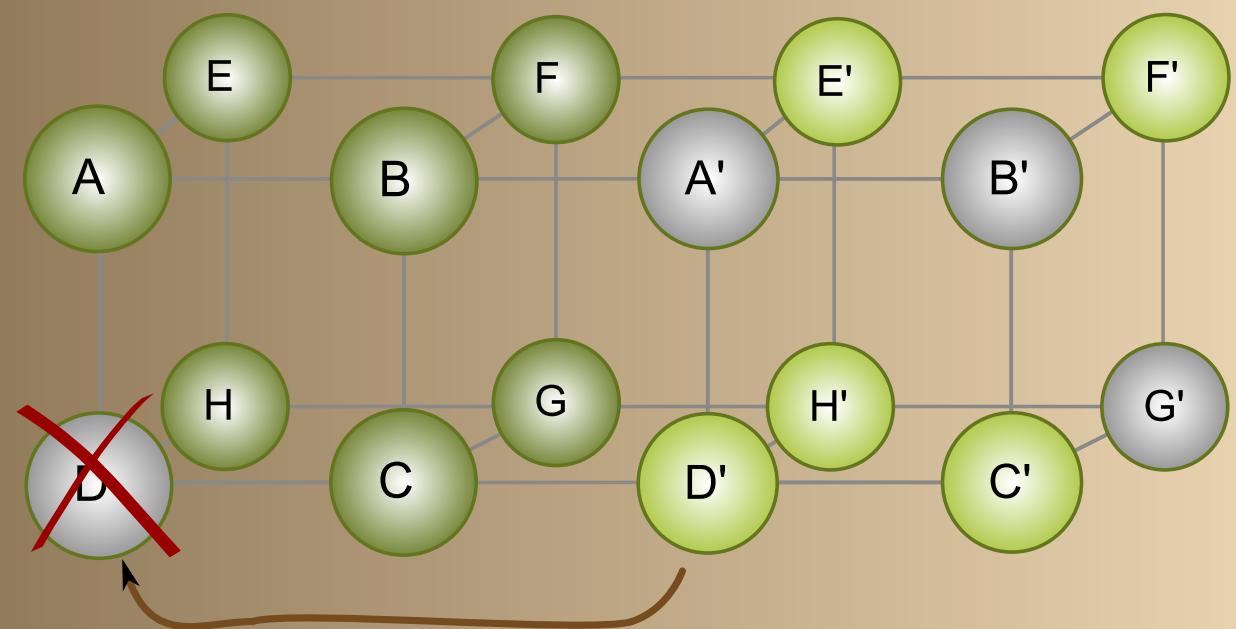
Status

Evaluation

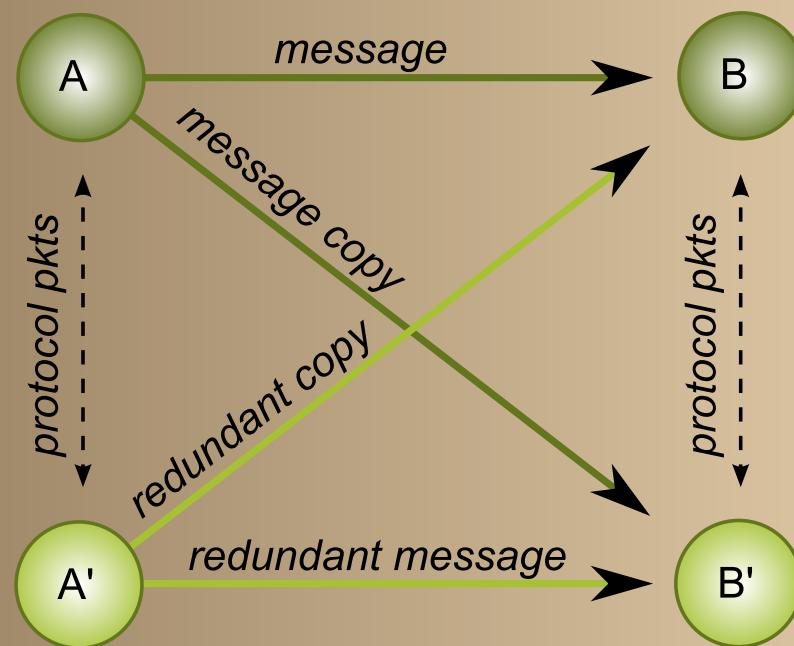
Analysis

Implications

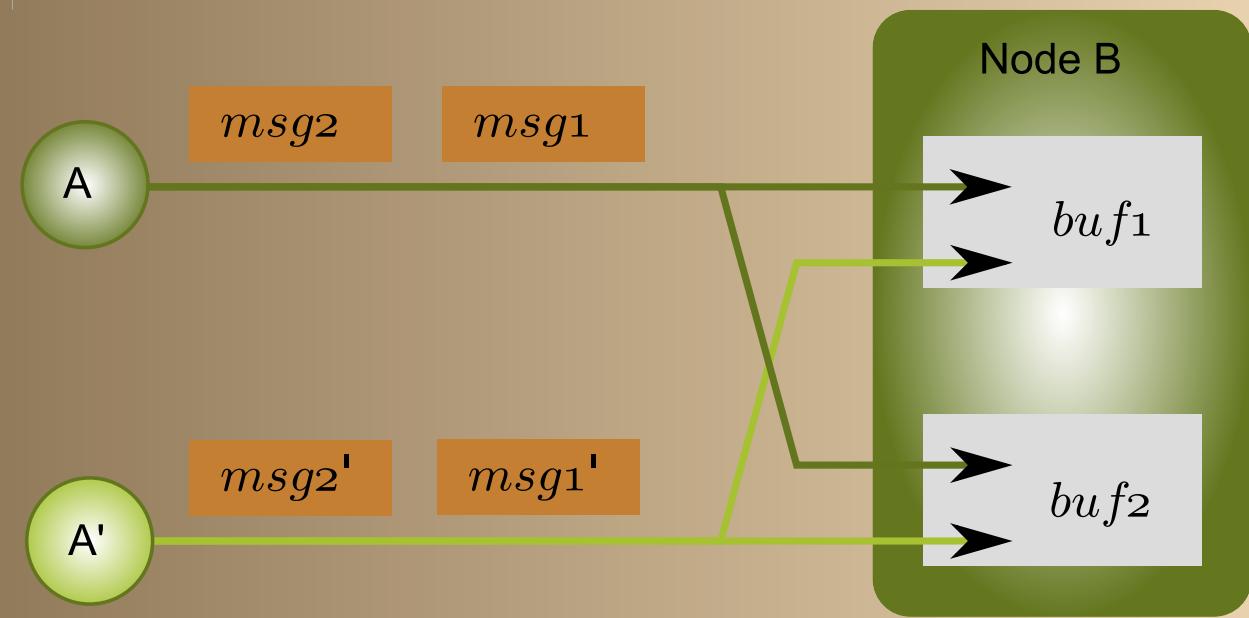
Summary



- Active and redundant node do same computation
- One node continues when the other fails
- MPI_Comm_rank() returns same value on both nodes
- Not each node needs to have a redundant partner



- Each message gets sent twice (4, if we count redundant nodes)
- Msg and redundant copy received into same buffer
- Redundant msg have unused tag bit set
- Protocol needed to coordinate receives and other MPI ops



- Active and redundant node must receive msg in same order
- MPI_ANY_SOURCE and MPI_ANY_TAG are problematic
- Redundant node maintains queue of posted receives if MPI_ANY_SOURCE
- Coordinate with active node to post specific receive

[Motivation](#)[Design](#)[Redundancy](#)[Basics](#)[Msg. order](#)[Other](#)[Status](#)[Evaluation](#)[Analysis](#)[Implications](#)[Summary](#)

- Redundant node must return same info for probe, test, and time functions
- Active node does operation and sends result to redundant node
- Collectives use redundant point-to-point
- *rMPI* re-implements almost all of MPI
 - ◆ *rMPI* uses MPICH (mostly) as a transport layer

[Motivation](#)[Design](#)[Redundancy](#)[Basics](#)[Msg. order](#)[Other](#)[Status](#)[Evaluation](#)[Analysis](#)[Implications](#)[Summary](#)

- Can't do MPI_ANY_TAG and MPI_ANY_SOURCE simultaneously
- Most major functions of MPI-2 implemented
- Are considering transactions to limit data volume
- Plan to open source it
- Few RAS features needed:
 - ◆ Notification of node availability
 - ◆ No Byzantine behavior (error correction protocol on network)
 - ◆ Messages to/from dead nodes must be consumed (no dead- or life-lock)

Motivation

Design

Evaluation

Bandwidth

Latency

Allreduce

CTH

SAGE

LAMMPS

HPCCG

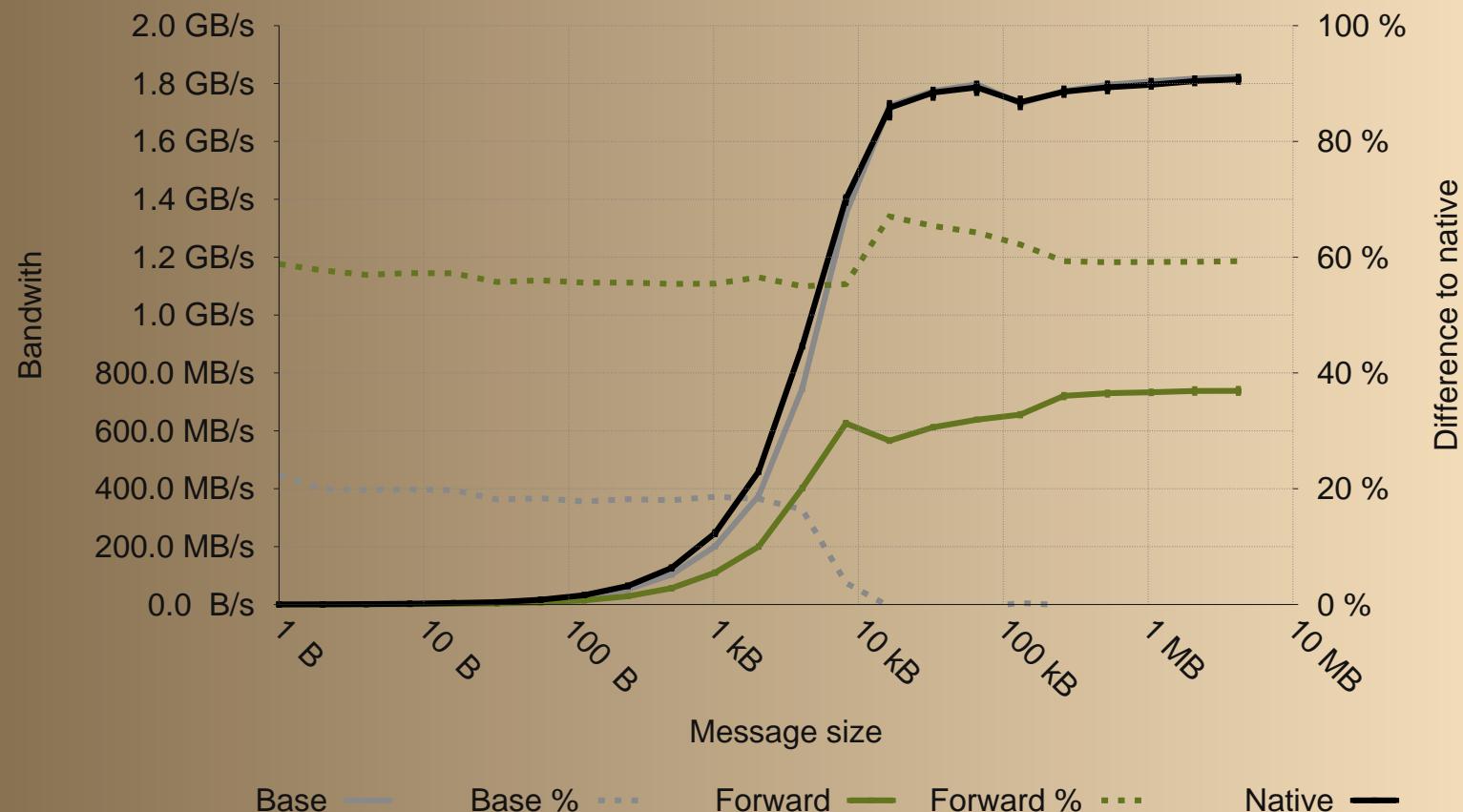
Analysis

Implications

Summary

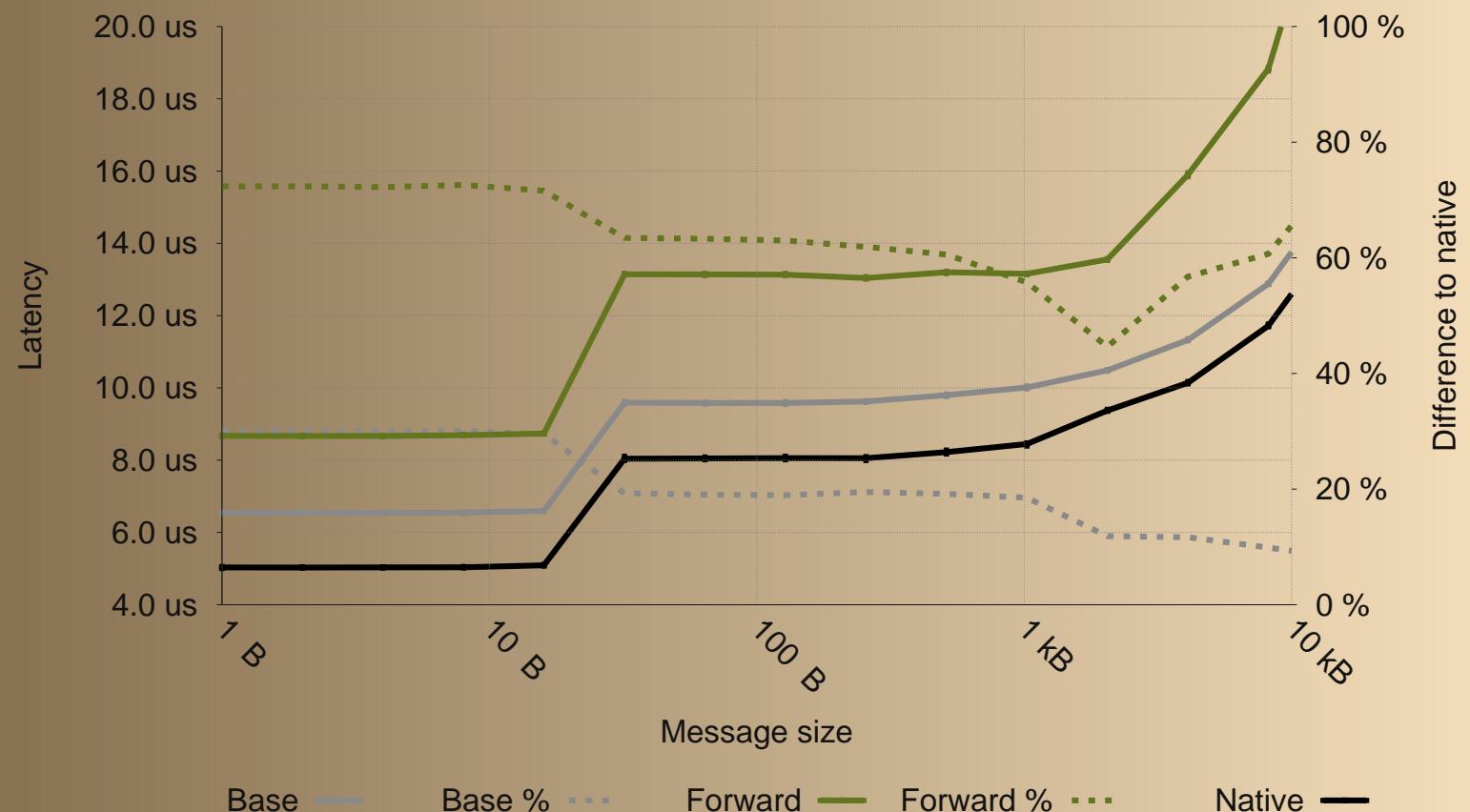
Evaluation

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



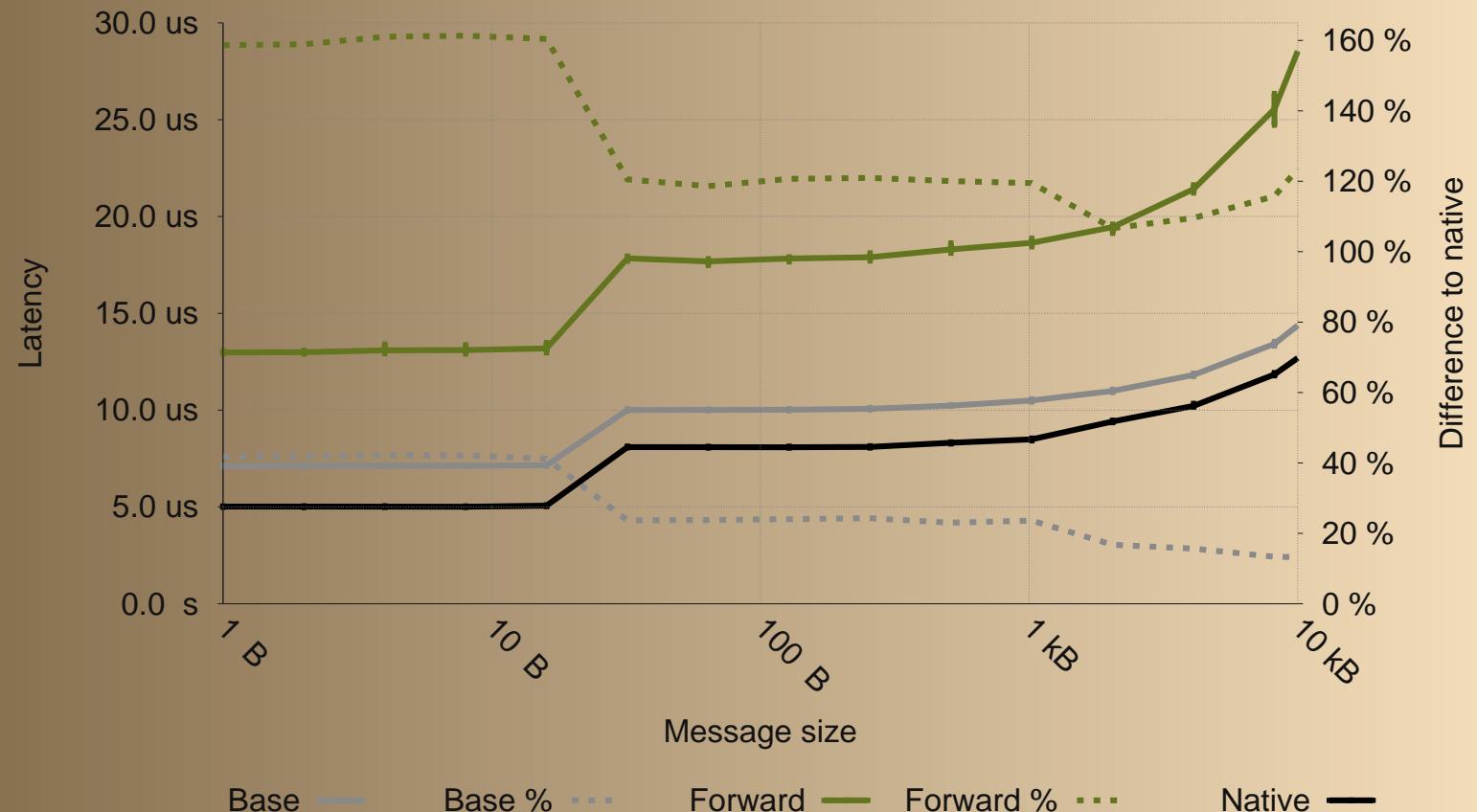
Native	Benchmark w/o rMPI
Base	rMPI, no redundancy
Forward	ABCD A'B'C'D'

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



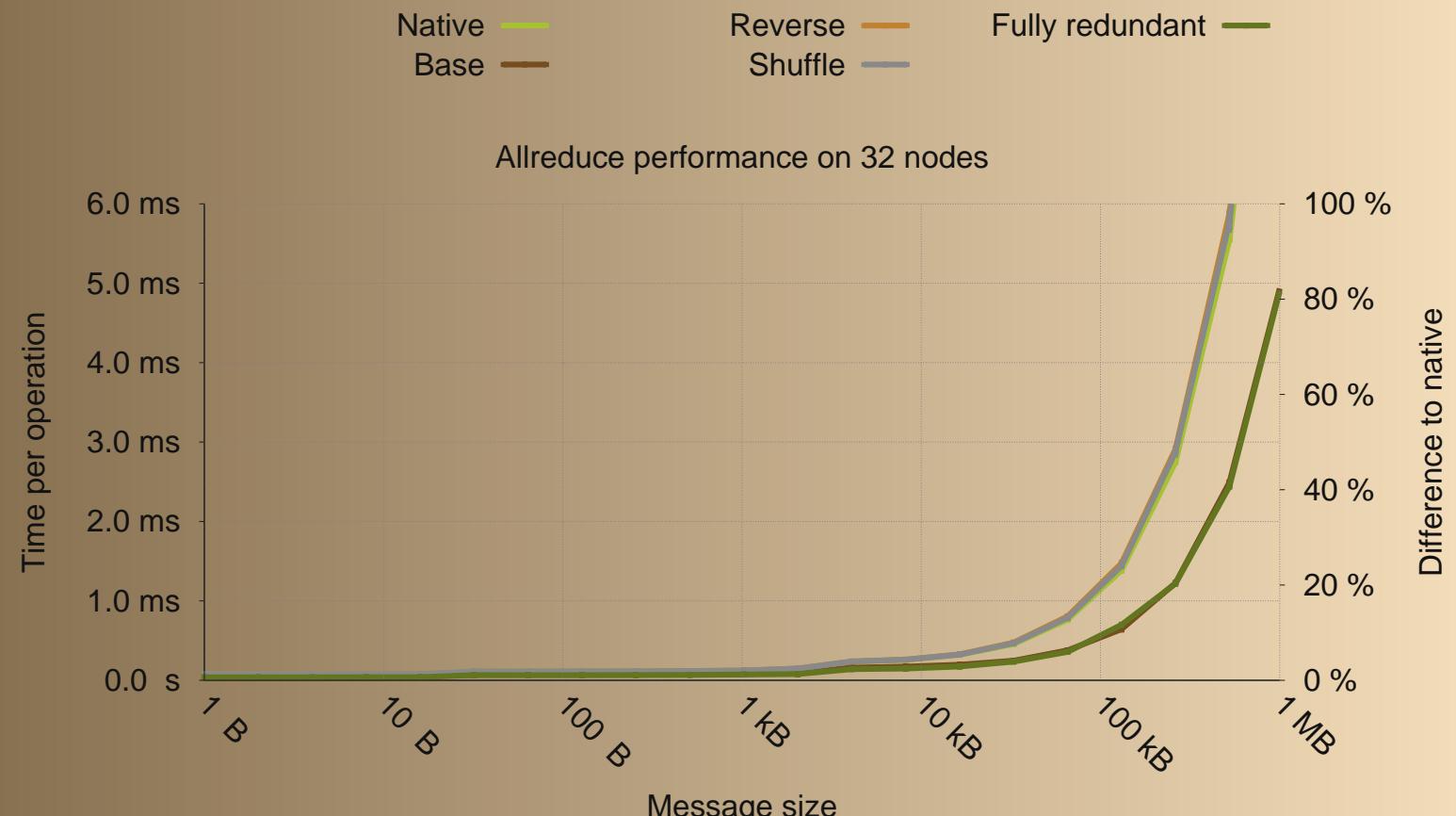
Native	Benchmark w/o rMPI
Base	rMPI, no redundancy
Forward	ABCD A'B'C'D'

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



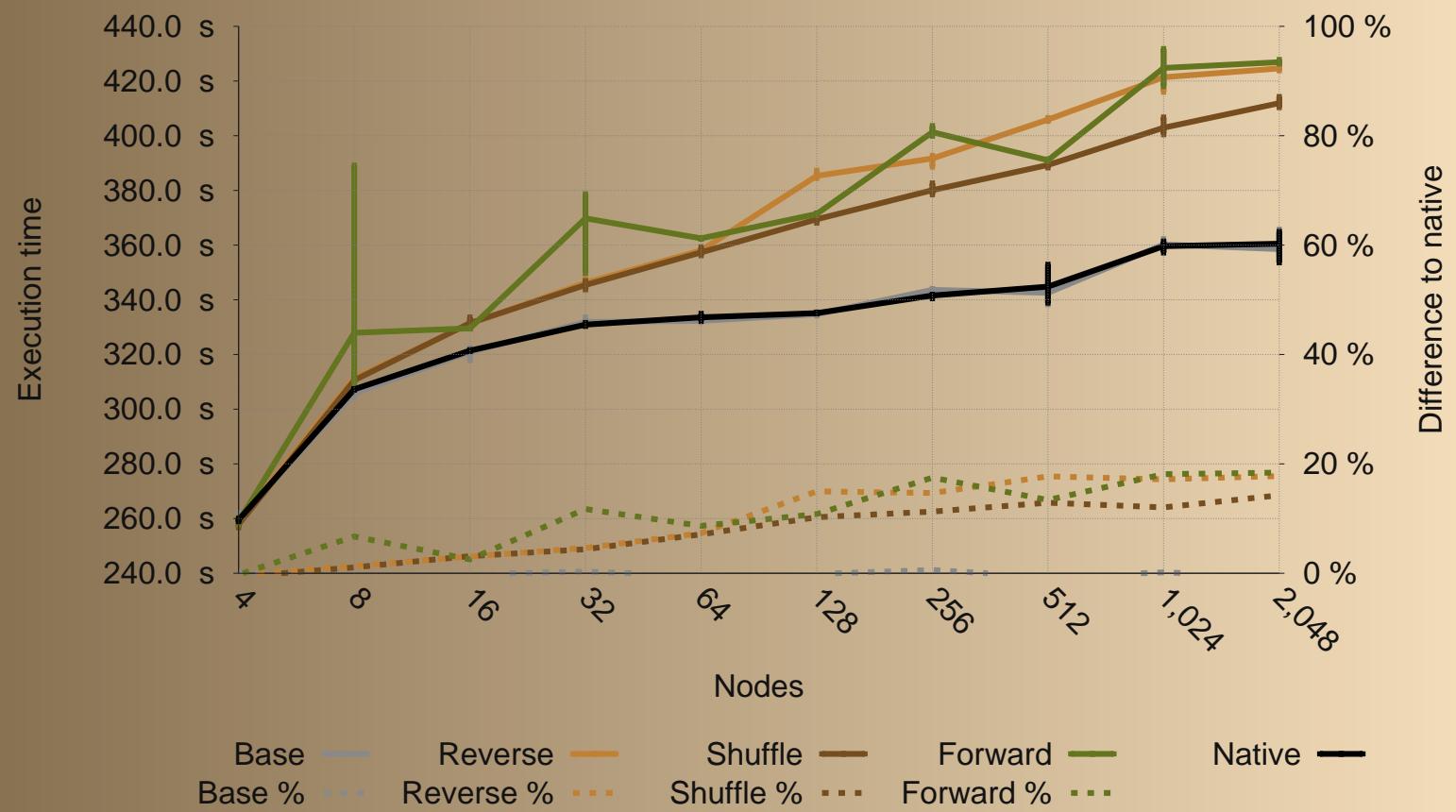
Native	Benchmark w/o rMPI
Base	rMPI, no redundancy
Forward	ABCD A'B'C'D'

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



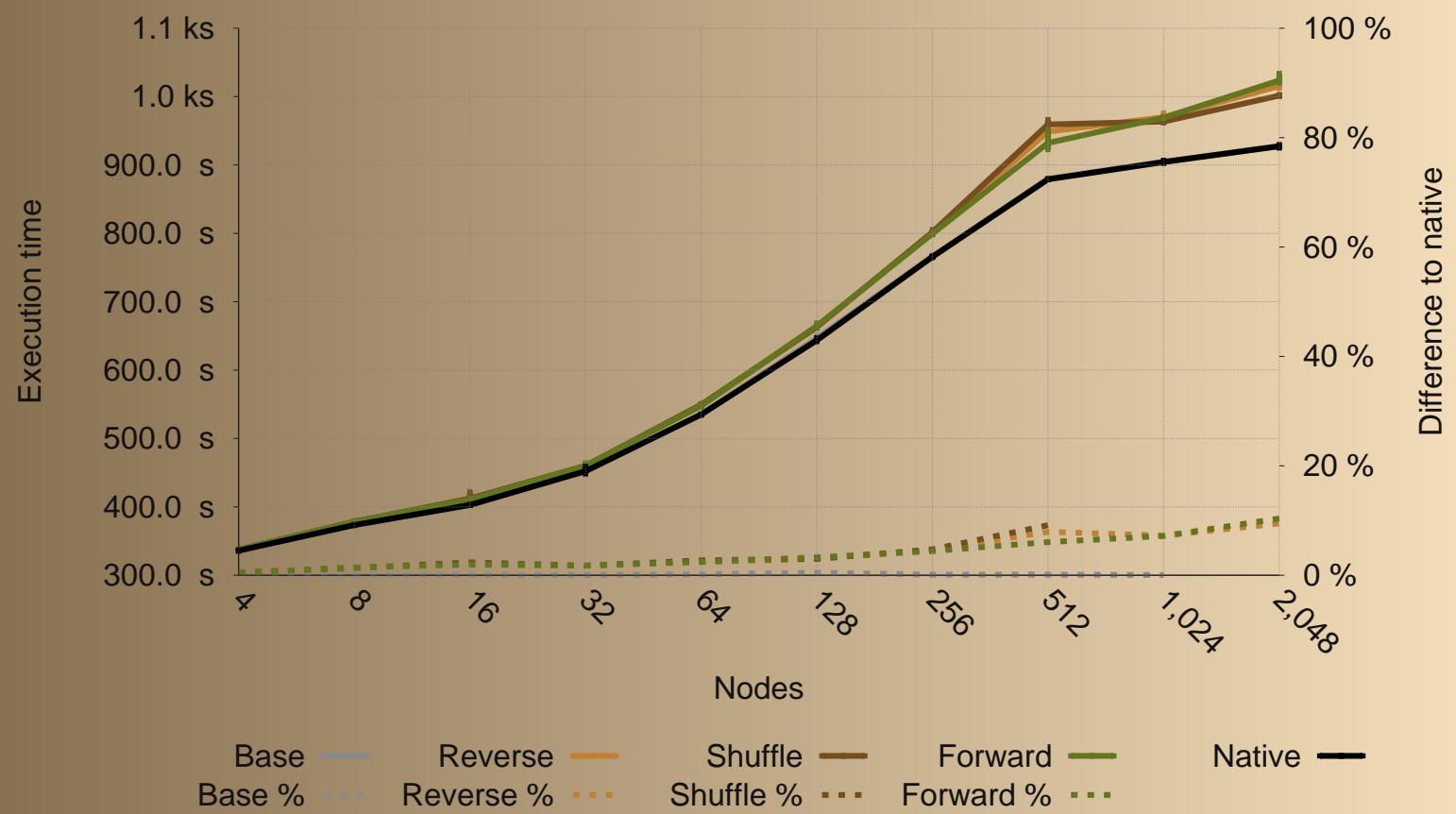
Native	Benchmark w/o <i>rMPI</i>	Reverse	ABCD D'C'B'A'
Base	<i>rMPI</i> , no redundancy	Shuffle	e.g., ABCD C'B'D'A'
Forward	ABCD A'B'C'D'		

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



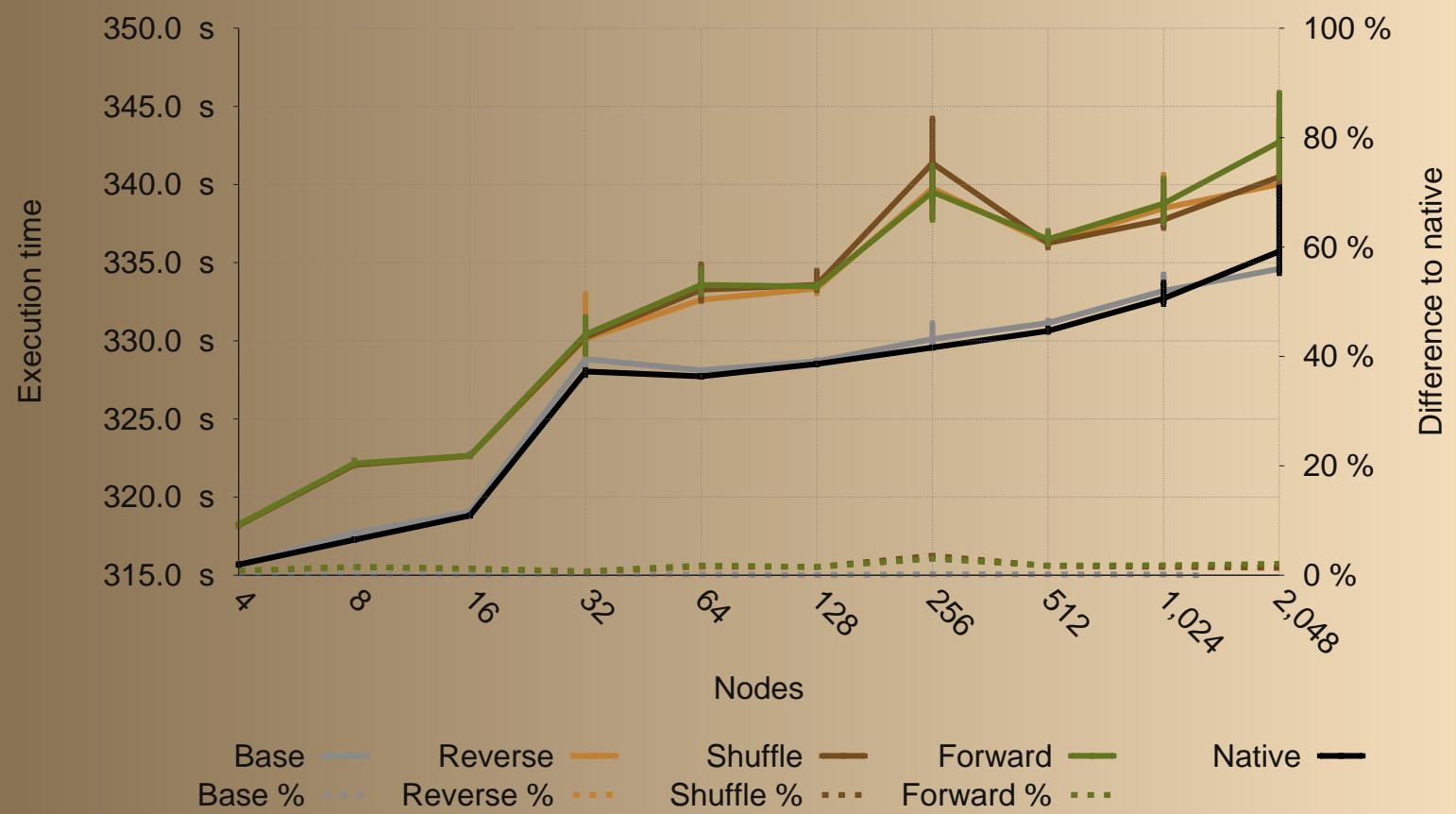
Native	Benchmark w/o rMPI		
Base	rMPI, no redundancy	Reverse	Shuffle
Forward	ABCD A'B'C'D'	ABCD D'C'B'A'	ABCD C'B'D'A'
		e.g., ABCD C'B'D'A'	

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



Native	Benchmark w/o <i>rMPI</i>	Reverse	ABCD D'C'B'A'
Base	<i>rMPI</i> , no redundancy	Shuffle	e.g., ABCD C'B'D'A'
Forward	ABCD A'B'C'D'		

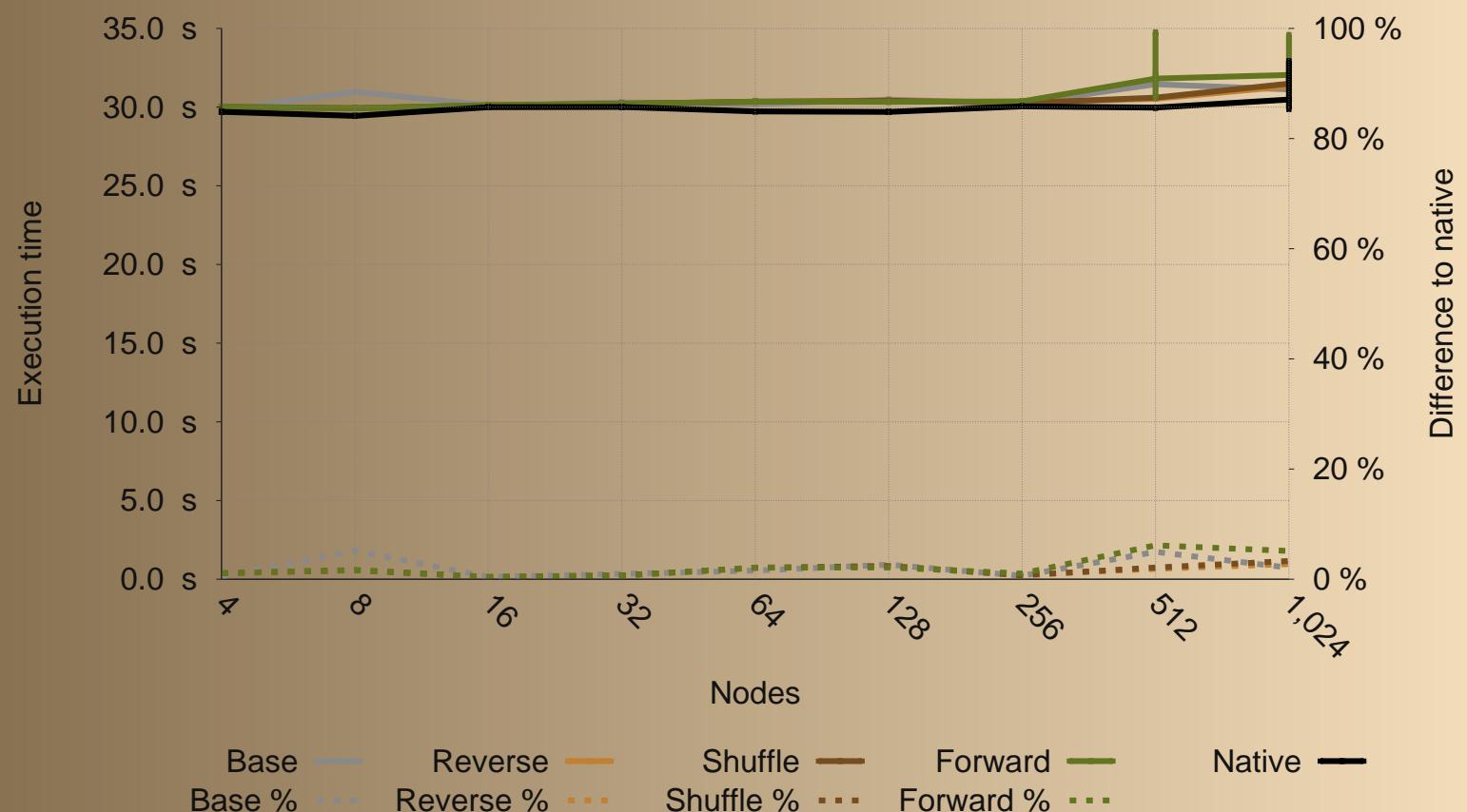
Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG
Analysis
Implications
Summary



Native	Benchmark w/o <i>rMPI</i>	Reverse	ABCD D'C'B'A'
Base	<i>rMPI</i> , no redundancy		
Forward	ABCD A'B'C'D'	Shuffle	e.g., ABCD C'B'D'A'

Motivation
Design
Evaluation
Bandwidth
Latency
Allreduce
CTH
SAGE
LAMMPS
HPCCG

Analysis
Implications
Summary



Native	Benchmark w/o <i>rMPI</i>	Reverse	ABCD D'C'B'A'
Base	<i>rMPI</i> , no redundancy	Shuffle	e.g., ABCD C'B'D'A'
Forward	ABCD A'B'C'D'	Forward	

Motivation

Design

Evaluation

Analysis

rMPI Model

Interrupts

Lifetime

Simulation

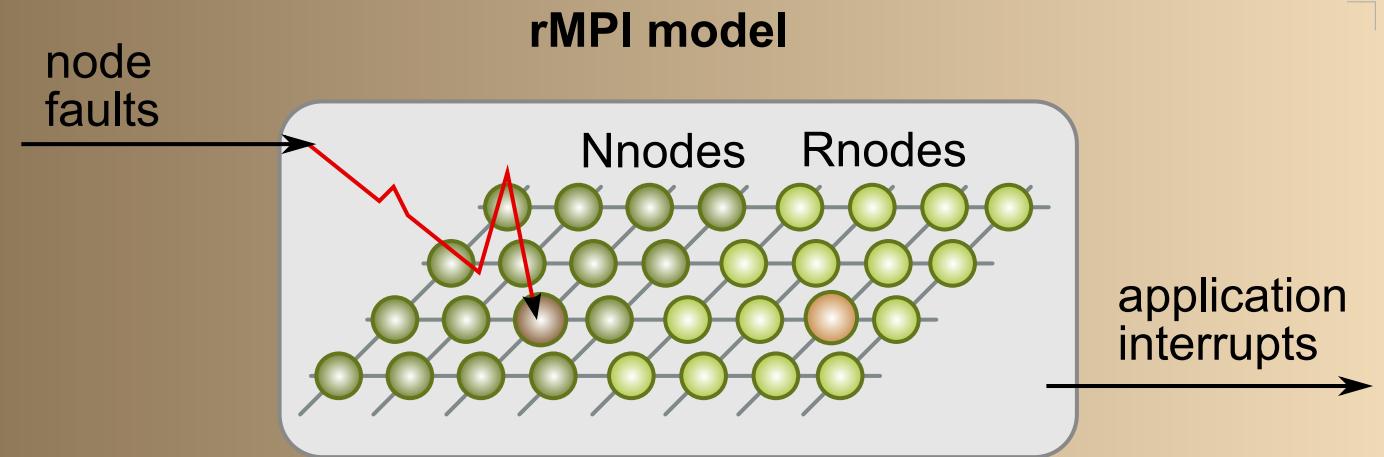
Behavior

Validation

Implications

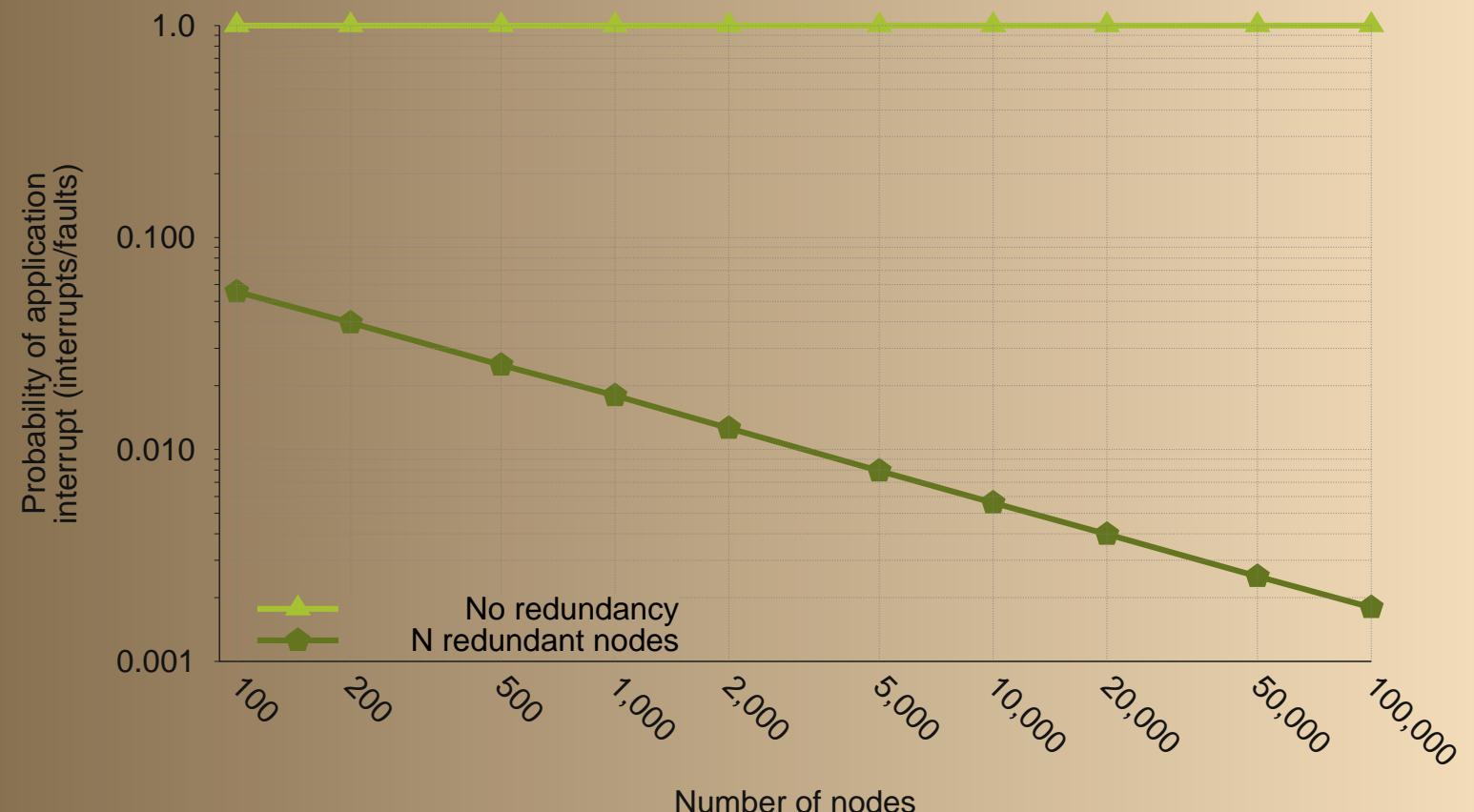
Summary

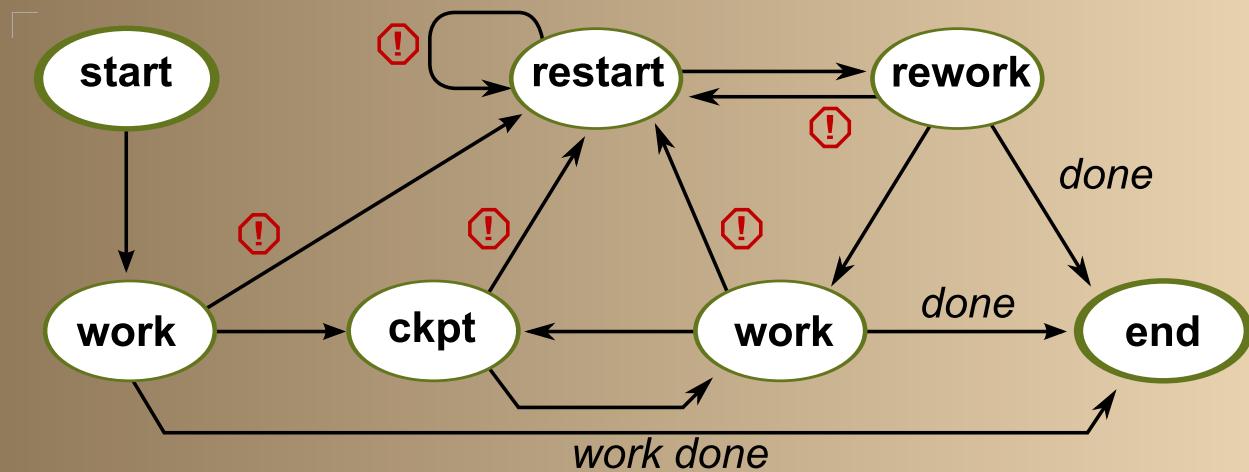
Analysis



- Acts as a filter
 - ◆ Input is a series of faults
 - ◆ Output is application interrupts that cause a restart

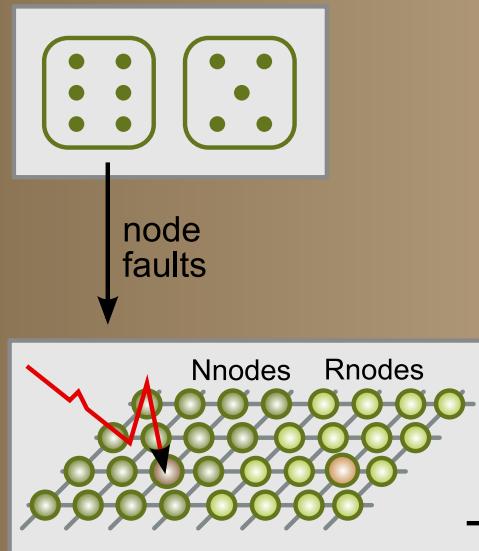
Motivation
Design
Evaluation
Analysis
rMPI Model
Interrupts
Lifetime
Simulation
Behavior
Validation
Implications
Summary





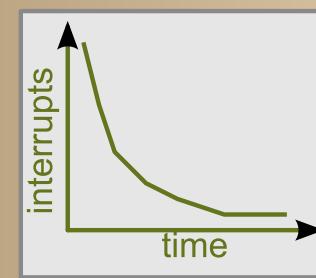
- Application needs to finish a set amount of work
- Do checkpoints and restart when necessary

Fault generator



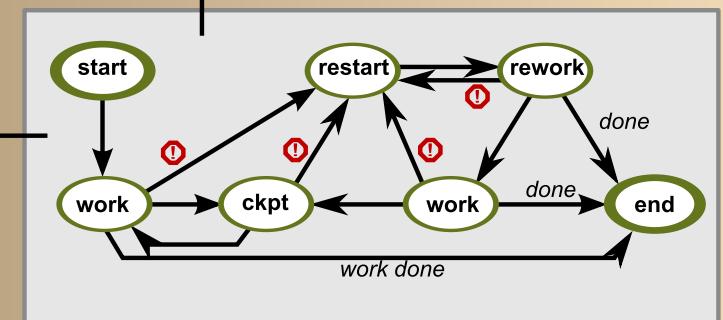
rMPI model

Analysis



statistical data

request next interrupt
application interrupts

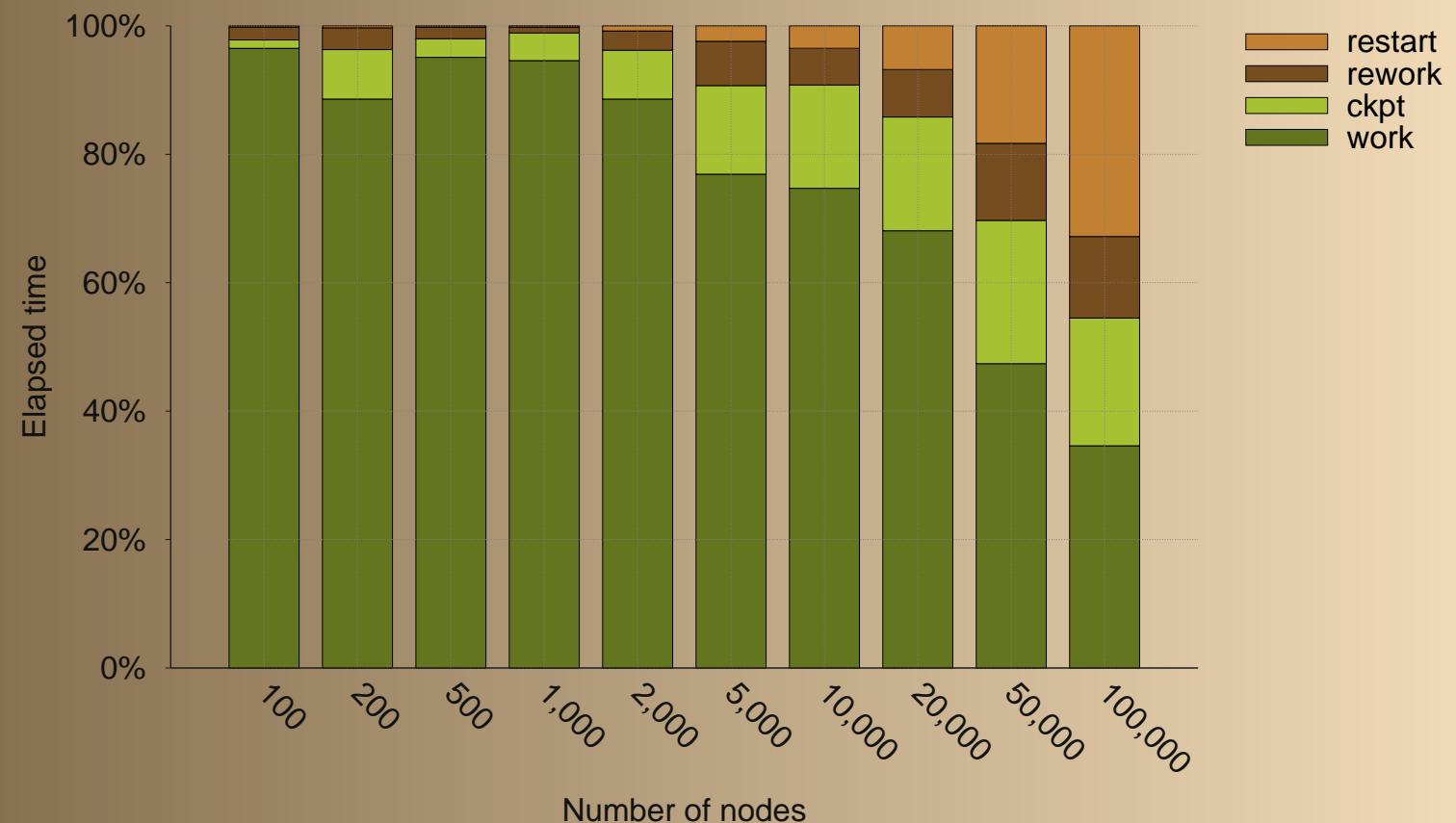


State machine

- Combine *rMPI* model with state machine
- Finish a set amount of work
- Do checkpoints and restart when necessary

Application Behavior no Redundancy

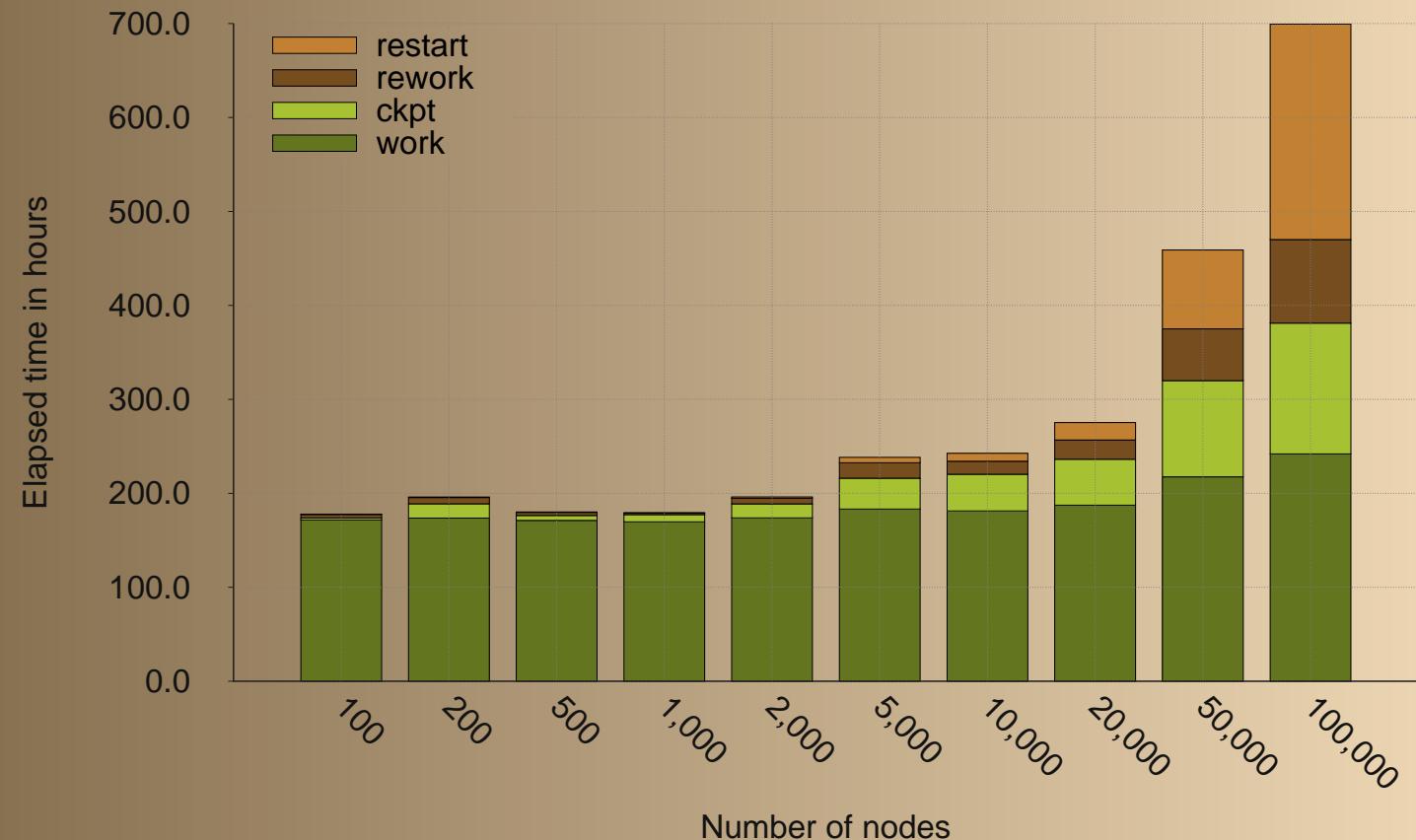
Motivation
Design
Evaluation
Analysis
rMPI Model
Interrupts
Lifetime
Simulation
Behavior
Validation
Implications
Summary



Normalized time spent in work, checkpoint, restart, and rework.
168h work, 5 min checkpoint, 10 min restart, 5-year node MTBF

Application Behavior no Redundancy

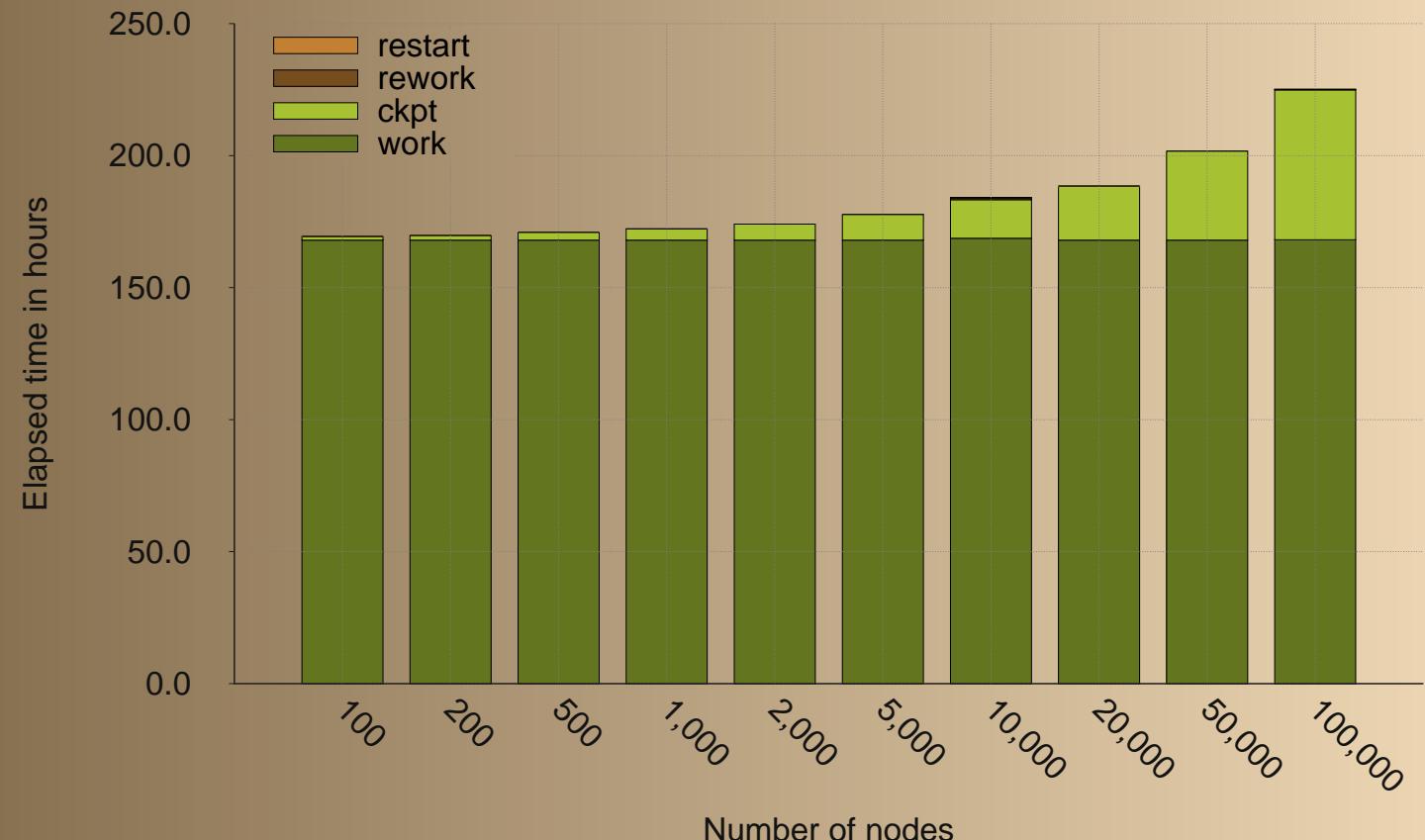
Motivation
Design
Evaluation
Analysis
rMPI Model
Interrupts
Lifetime
Simulation
Behavior
Validation
Implications
Summary



Total time spent in work, checkpoint, restart, and rework.
168h work, 5 min checkpoint, 10 min restart, 5-year node MTBF

Application Behavior with Redundancy

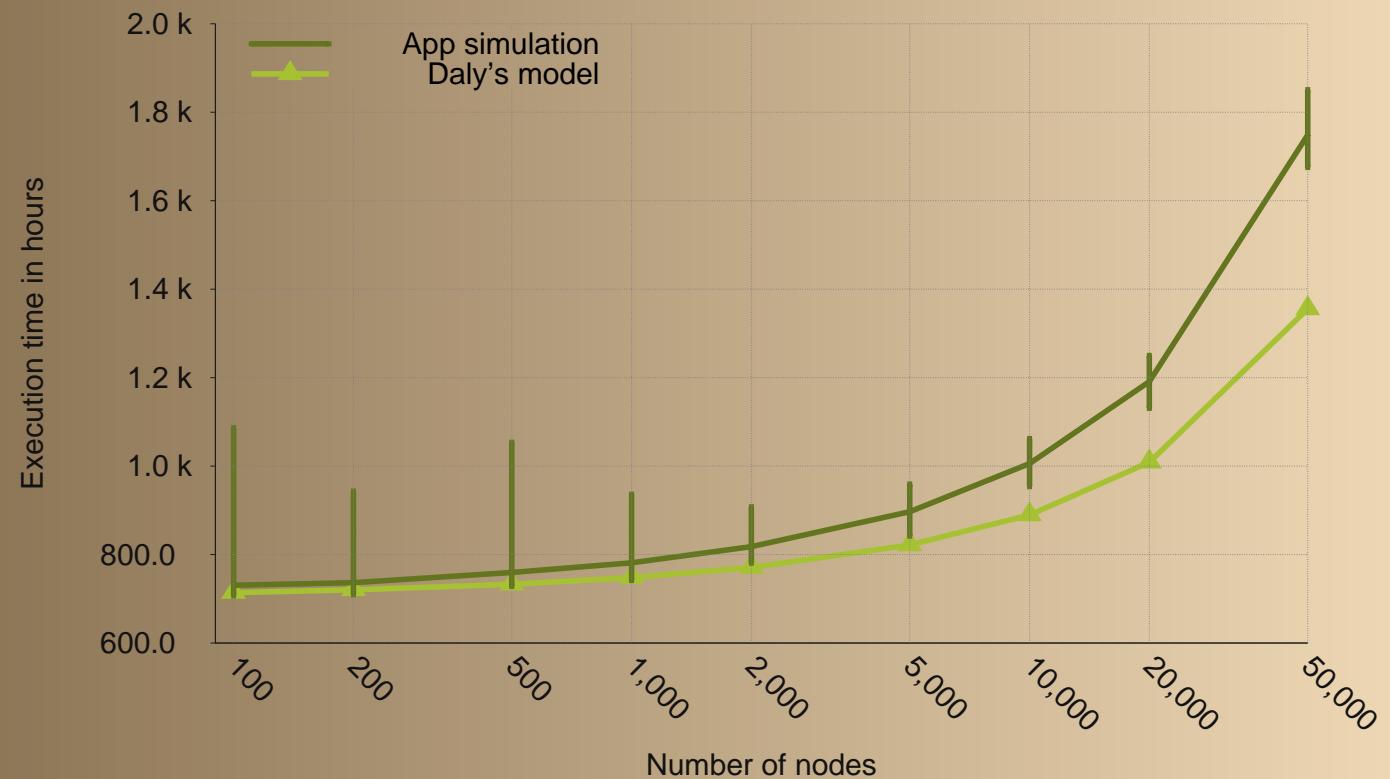
Motivation
Design
Evaluation
Analysis
rMPI Model
Interrupts
Lifetime
Simulation
Behavior
Validation
Implications
Summary



Normalized time spent in work, checkpoint, restart, and rework.
168h work, 5 min checkpoint, 10 min restart, 5-year node MTBF

Daly's equation from *A higher order estimate of the optimum checkpoint interval for restart dumps.*

$$T_w(\tau) = \Theta e^{\frac{R}{\Theta}} (e^{\frac{\tau+\delta}{\Theta}} - 1) \frac{T_s}{\tau} \quad \text{for } \delta \ll T_s$$



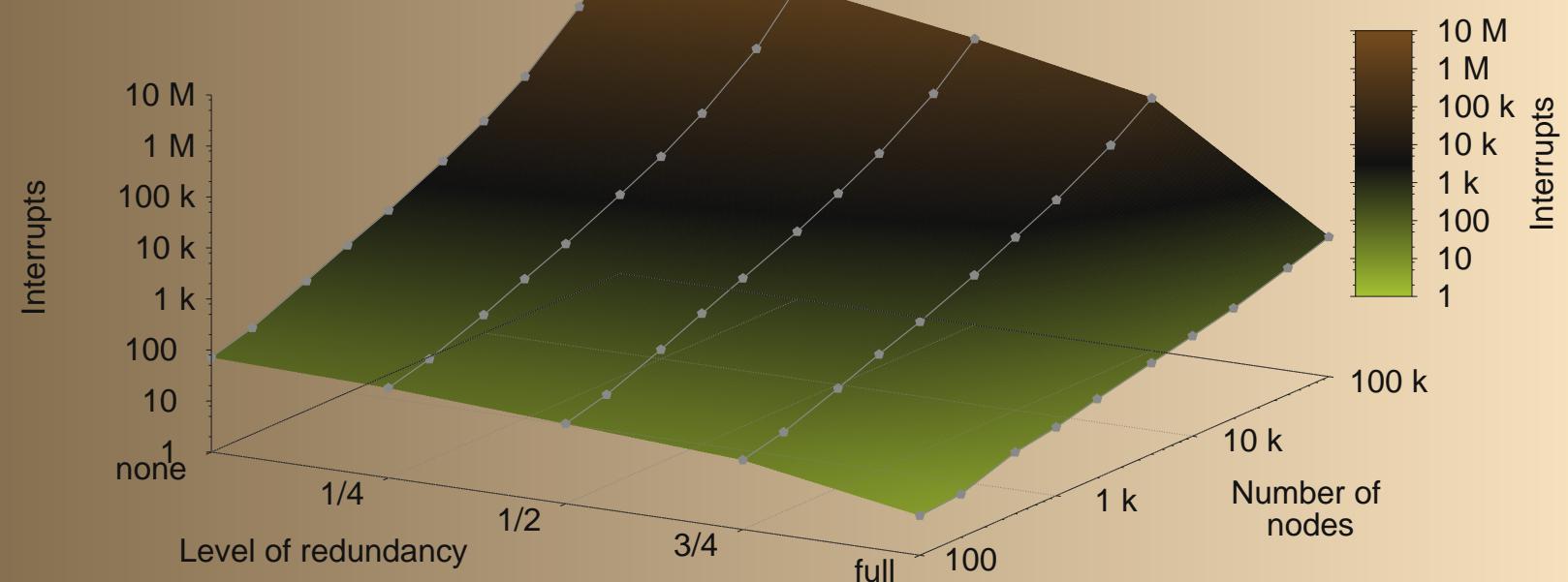
Motivation
Design
Evaluation
Analysis
Implications
Summary

Implications and Trade-Offs

$$R = \frac{T_w + T_{o(\text{none})}}{T_w + T_{o(\text{full})} + T_{r\text{MPI}}} \quad (1)$$

- T_w : amount of work
- $T_{o(\text{none})}$: checkpoint, restart, rework overhead without redundant nodes
- $T_{o(\text{full})}$: overhead with redundant nodes
- If $R > 2$ then using redundant nodes makes sense

Motivation
Design
Evaluation
Analysis
Implications
Summary



5,000h work, 5 min checkpoint, 10 min restart, 1-year node MTBF

Motivation

Design

Evaluation

Analysis

Implications

Summary

Summary

- Redundant MPI library can be done at user level
- Overhead for applications is not significant for most applications
- Application restart simulator allows modeling of
 - ◆ varying node counts
 - ◆ node MTBF
 - ◆ level of redundancy
 - ◆ failure distribution function (exponential, gamma, weibull)
 - ◆ amount of work, checkpoint and restart times
- Model helps determine when redundancy pays off